

[> home](#) [> about](#) [> feedback](#) [> login](#)

US Patent & Trademark Office

Citation

ACM SIGCOMM Computer Communication Review [>archive](#)Volume 23 , Issue 1 (January 1993) [>toc](#)**Extending the IP internet through address reuse**

Authors

Paul F. Tsuchiya
Tony Eng



Publisher

ACM Press New York, NY, USA

Pages: 16 - 33 Periodical-Issue-Article

Year of Publication: 1993

ISSN:0146-4833

 <http://doi.acm.org/10.1145/173942.173944> (Use this link to Bookmark this page)[> index terms](#)[> Discuss](#)[> Similar](#)[> Review this Article](#) [Save to Binder](#)[> BibTex
Format](#)

[↑ INDEX TERMS](#)**Primary Classification:**

C. Computer Systems Organization



C.2 COMPUTER-COMMUNICATION NETWORKS



C.2.1 Network Architecture and Design

Extending the IP Internet Through Address Reuse

Paul F. Tsuchiya, Bellcore, tsuchiya@thumper.bellcore.com

Tony Eng, MIT, tleng@athena.mit.edu

Abstract

The two most compelling problems facing the IP Internet are IP address depletion and scaling in routing. This paper discusses the characteristics of one of the proposed solutions—address reuse. The solution is to place Network Address Translators (Nat) at the borders of stub domains. Each Nat box has a small pool of globally unique IP addresses that are dynamically assigned to IP flows going through Nat. The dynamic assignment is coordinated with Domain Name Server operation. The IP addresses inside the stub domain are not globally unique—they are reused in other domains, thus solving the address depletion problem. The pool of IP addresses in Nat is from a subnet administered by the regional backbone, thus solving the scaling problem. The main advantage of Nat is that it can be installed without changes to any existing systems, although FTP will fail in some but not all cases. This paper presents a preliminary design for Nat, and discusses its pros and cons.

1.0 Introduction

The two most compelling problems facing the IP Internet are IP address depletion and scaling in routing. Numerous solutions have been proposed, such as increasing the size of the IP address, using completely flat IP addresses, changing the structure of IP addresses, and even abandoning IP and switching to OSI [Ch]. Unfortunately, all of these solutions require changes to routers, hosts, or both.

Among the proposed solutions is the notion of address reuse. This solution takes advantage of the fact that a very small percentage of hosts in a stub domain¹ are communicating outside of the domain at any given time. Indeed, many (if not most) hosts never communicate outside of their stub domain. Because of this, IP addresses inside a stub domain, that usually need not be globally unique or known externally, can be dynamically translated into a small pool of IP addresses that are globally unique when outside communications is required.

1. A stub domain is a domain, such as a corporate network, that only handles traffic originated by or destined to hosts in the domain.

This solution has the disadvantage of taking away the end-to-end significance of an IP address, and making up for it with increased state in the network. There are various work-arounds that minimize the potential pitfalls of this. Indeed, connection-oriented protocols are essentially doing a kind of address reuse at every hop.

The huge advantage of this approach is that it can be installed incrementally, without changes to either hosts or routers². Depending on how Nat is implemented, changes to the Domain Name System (DNS) server in the stub domain may be required. This solution can be implemented and experimented with quickly. If nothing else, this solution can serve to provide temporarily relief while other, more complex and far-reaching solutions are worked out.

2.0 Overview of Nat

The design presented in this paper is called Nat, for Network Address Translator. Nat is a box or router function that can be configured as shown in figure 1. The upper configurations require no host or router modifications. The lower configuration requires a modification to the stub border router.

Nat's basic operation is as follows. The addresses inside a stub domain can be reused by any other stub domain. At each exit point between a stub domain and backbone, Nat is installed. Each Nat box is assigned a small pool of globally unique IP addresses (each Nat has a separate pool). These IP addresses are dynamically assigned to IP "flows" going through Nat.

For instance, in the example of figure 2, both stubs A and B internally use class A network number 42.0.0.0. Stub A's Nat is assigned the (subnetted) class B subnet number 128.76.29.0, and Stub B's Nat is assigned the class B subnet number 128.76.28.0³. The class B subnets are globally unique—no other Nat boxes can use them.

When stub A host 42.33.96.5 wishes to exchange packets with stub B host 42.81.13.22 ("al.nxb.com"), it sends a Domain Name System (DNS) query to the DNS in stub B (1). DNS knows that the internal address for "al.nxb.com" is 42.81.13.22, but let's assume that there is no external address assigned for "al.nxb.com". DNS would then send a query to Nat asking to have an address assigned (2). Nat finds an

2. A few unusual applications may require changes, and hosts that communicate outside their domain to hosts that do not have permanent assignments must use DNS.

3. If a backbone assigned subnetted portions of a class B, then it could represent multiple stubs by advertising a single class B address externally.

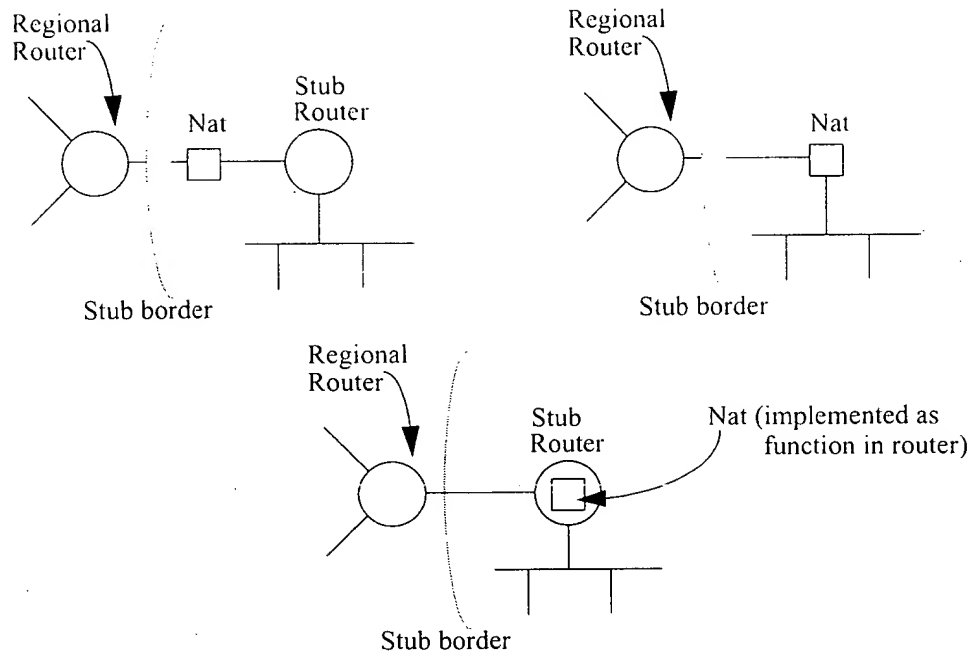


Figure 1: Nat Configurations

unassigned address in its pool, 128.76.28.4, and returns it to DNS (3), that in turn uses this address to answer the original DNS query (4). Alternatively, the DNS box could not query the Nat box, but rather send the answer (42.81.13.22) back towards stub A's DNS. The Nat box would intercept the DNS answer, assign a temporary address (128.76.28.4), and modify the DNS answer to return 128.76.28.4 to the DNS server in stub A. This latter approach requires no changes to DNS, but constrains the configuration of Nat boxes to one clique (see section 3.4). This constraint will not affect most domains.

42.33.96.5 then sends a packet to "al.nxb.com" with destination address 128.76.28.4. When this packet reaches stub A's Nat, it assigns an externally unique address (128.76.29.7) from the pool to 42.33.96.5, and translates the source address of the IP header with the new address⁴. This packet is routed through the backbones to the stub B Nat, which translates the destination address of the IP header to be the internally known address, and the packet is sent to "al.nxb.com". Likewise, IP packets on the return path go through similar address translations.

4. Note that both the IP and TCP checksums must be modified. This does not require a complete recalculation—only an incremental recalculation.

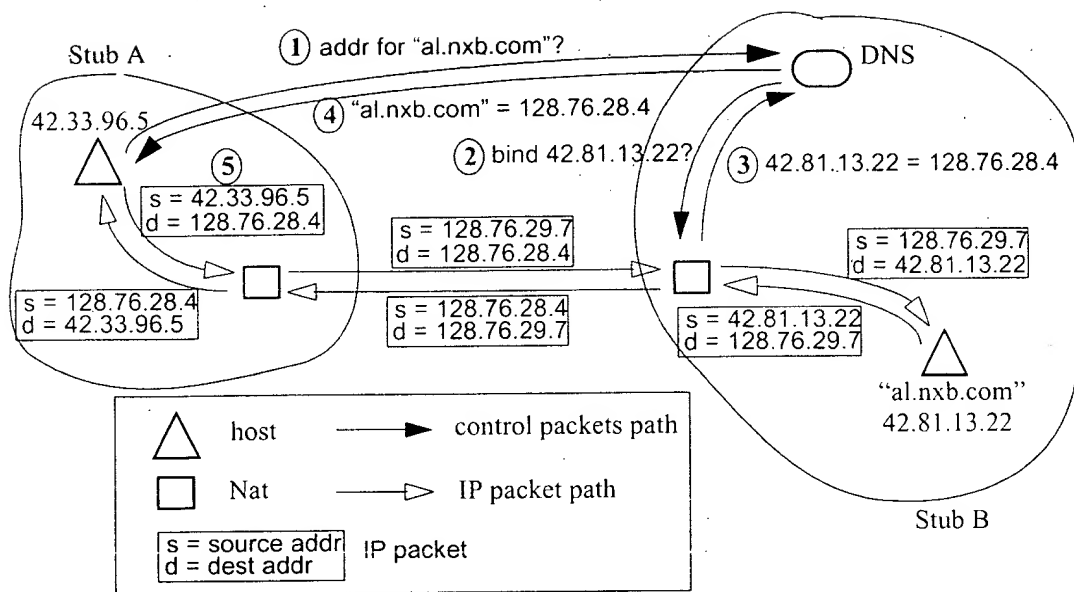


Figure 2: Basic Nat Operation

Notice that this requires no changes to hosts or routers. For instance, as far as the stub A host is concerned, 128.76.28.4 is the address used by "al.nxb.com". The address translations are completely transparent. Depending on whether or not the Nat box intercepts DNS packets, only the Domain Name Server may require modification.

Of course, this is just a simple example. There are numerous issues to be explored. In the next section, we discuss various aspects of Nat.

3.0 Various Aspects of Nat

3.1 Address Spaces

Partitioning of Reusable and Non-reusable Addresses

For Nat to operate properly, it is necessary to partition the IP address space into two parts—the reusable addresses used internal to stub domains, and the globally unique addresses. We call the reusable address *local addresses*, and the globally unique addresses *global addresses*. Any given address must either be a local address or a global address. There is no overlap.

The problem with overlap is the following. Say a host in stub A wished to send packets to a host in stub B, but the global addresses of stub B overlapped the local addressees of stub A. In this case, the routers in stub A would not be able to distinguish the global address of stub B from its own local addresses, and would not know whether to route the packets internally or externally.

Initial Assignment of Local and Global Addresses

Theoretically all stubs could use the same class A address locally. However, existing stubs already have unique addresses assigned internally. It is difficult and takes time to change all addresses in a stub. Therefore, at least initially, existing address assignments should be defined as local addresses. A block of unassigned class B addresses should be defined as global addresses. These would be assigned to Nat boxes (on a subnetted basis). A single class A address should also be defined as local. This class A would be given to new stubs, who would be expected to install Nat when they connect to the IP Internet. Over time, existing stubs should install Nat and transition their existing address to the class A address. Once the transition was complete, the stub could give back its old addresses, which would then become global.

Scaling

One assignment strategy for global addresses goes as follows. Each regional (bottom level backbone) would be assigned a global class B address. The regional would then subnet the class B address among Nat boxes that connected to it. If, for instance, each Nat box required an average of 250 global addresses in its pool (not an unreasonable estimate⁵), then the class B address could be subnetted among 250 or so Nat boxes. Even if most stubs had two connections to the regional (thus requiring two Nat boxes, and two pools of addresses), 125 stubs could be subnetted into one class B address. The regional could then advertise one class B address to other backbones, rather than 125 separate addresses, as it does now. This would shrink current routing tables from several thousand entries to tens of entries.

3.2 Address Assignments

Permanent Address Assignments

Not all address assignments made by Nat should be dynamic. Any hosts that communicate outside the stub frequently should be given permanent assignments. Such hosts include DNS servers, email distributors, and anonymous FTP repository hosts. Indeed, stubs that implement security by limiting outside

5. During one week of measurements of external TCP connections at Lawrence Berkeley Laboratory [Pa], the largest number of active simultaneous TCP connections registered was 57.

communications to a small number of “secure” hosts do not need dynamic assignment at all. Of course, the DNS servers **must** have permanent assignments, because it is through the DNS servers that dynamically assigned global addresses become known.

The global address pool (or just pool), then, is partitioned into two parts, the static pool and the dynamic pool.

Choosing an Assignment

When Nat intercepts a DNS response (or receives a request from DNS for an assignment⁶), or when a packet that does not have an assignment arrives from the stub destined for the backbone, Nat must choose an assignment from the dynamic pool. The task of choosing an assignment from the pool can be tricky. The goal is to 1) minimize prematurely destroying assignments while 2) maximizing address utilization and 3) minimizing complexity. The primary goal is the first one, provided of course that we can make significant efficiency gains in address utilization over current practice.

The simplest algorithm for choosing an assignment is to choose the address that has not been in active use the longest. To do this, Nat would save the time that the last packet was seen for each of its assignments. When a new assignment is needed, Nat chooses the one with the oldest time. We call this algorithm *IP-based*, because it bases its activity estimate purely on the timing of individual IP packets.

The IP-based algorithm assumes (incorrectly, of course) that the usage behavior of all assignments is the same. For instance, it would choose the assignment for a TCP connection that has been idle for two minutes but never issued a close, over a TCP connection that issued a close one minute ago and has been idle ever since. Clearly, it would be preferable to destroy the one-minute old closed connection over the two-minute old open connection⁷. Since the IP-based algorithm does not monitor connection status, it cannot know this.

If applications where a valid assignment may remain idle for long periods of time exist, then the IP-based algorithm must maintain a large pool to insure that the oldest assignment is not still active. (If a limited

6. Throughout this paper, reference may be made to either the intercept method of modifying DNS packets, or explicit packet exchange between DNS boxes and Nat boxes, but not both. Unless otherwise stated, it should be understood that either method is possible.

7. TCP does not have a keep-alive function, so a TCP connection can remain open without exchanging packets indefinitely.

number of hosts have these kinds of applications, then a possible solution would be to give these hosts permanent assignments.)

An alternative approach would be an algorithm that monitors connection status and partitions assignments into two classes—those that have seen an end-of-connection indication and those that have not. The ended-connection assignments would become available for re-assignment after a relatively short idle time, say one minute. The open assignments would become available for re-assignment after a long idle time, perhaps several hours. One way to implement this would be to consider the working idle time of an ended-connection association to be x times the actual idle time, and then compare all idle times directly. We call this the *connection-based* algorithm.⁸

It is impossible to know which assignment algorithm to use without knowing how applications behave. Since Nat will initially be experimental, both algorithms can be experimented with. Since each Nat box can independently choose its own algorithm, multiple algorithms can be experimented with simultaneously.

Incorrect Assignments

The reuse of global addresses by Nat is not explicitly coordinated with the use of addresses by hosts. Through DNS packets, Nat can inform a host of an assignment. However, Nat cannot inform hosts when an address has been re-assigned. As a result, it is entirely possible for a host A to use an address that it thinks is for host B, but that in fact Nat has reassigned to host C.

Note that even without address reuse, there always exists the possibility that a host uses the wrong address. For instance, DNS can be incorrectly configured so that the IP address it returns does not belong to the host named in the query. Usually, the application protocol or the human user will discover the error. Nat, however, increases the probability of misdelivering packets.

Notice that a protocol such as X.25 is basically an address reuse scheme similar in some respects to Nat. In this case, the X.121 address corresponds to the Domain Name, and the Virtual Circuit Identifier (VCI) corresponds to the assigned IP address. X.25 for all practical purposes does not have the mis-addressed packet problem like Nat does, because the assignment of VCIs is explicitly coordinated by all components on the path, including hosts. It follows, then, that the problem of mis-addressed packets in Nat is not a problem with address reuse per se, but a problem of the style of implementation resulting from the decision to keep Nat transparent to hosts.

8. Connection-less applications can only use the IP-based algorithm

All cases of mis-addressed packets are the result of hosts (or DNS) caching addresses longer than Nat. This comes about because Nat, through the modification of DNS packets, can give a host an address to use, but cannot later remove it. Nat should therefore always set the Time-to-Live (TTL) in DNS packets (or responses to DNS queries for an assignment) to a value slightly smaller than the minimum time that Nat will maintain an assignment [Lo]. This of course does not prevent a host from caching the assignment for a longer period of time. For instance, the host may start up an application that runs for a long time, sends very occasional UDP packets, always using the IP address that it was originally invoked with. A small TTL only prevents a DNS server from storing the query beyond this time, thereby preventing it from giving the assignment to another host after it has expired.

Another style of implementation, that requires changes to DNS and the involvement of DNS at both source and destination, but that reduces the probability of mis-addressed packets, is briefly described as follows. Instead of distinguishing assignments by only stub-side addresses, Nat distinguishes on both sides. In other words, if a packet does not have an expected source *and* destination address, it is dropped.

When a host H_a wants to send packets, it queries its local DNS server D_a . Rather than immediately send a query to the destination DNS server, D_a queries the Nat box N_a in its domain and gets an assignment A_a . This assignment is then conveyed in the query to the DNS server D_b on the destination side. D_b then conveys the address A_a to the destination Nat box N_b , which associates the source global address A_a with the destination global address A_b that it assigns. When D_b answers the query, D_a informs N_a of the address A_b , so now both Nat boxes have both A_a and A_b associated with their assignments. As a result, hosts that previously used either A_a or A_b will now not be able to use them, and packets will not be mis-delivered (just dropped). To successfully communicate, these hosts will have to again query their DNS server.

If a third host H_c wants to send packets to host H_b with global address A_c , DNS server D_b will need to inform N_b of the new address A_c so that N_b will now associate both A_a and A_c with A_b . Note that if only one side has implemented Nat, the DNS server on the Nat side will need to query the DNS server on the non-Nat side to learn the expected IP address of the other side. This extra query is needed because DNS queries normally only carry the domain name of the query originator, not the IP address.

This two-sided design involves more overhead and complexity than the one-sided design, and may turn out not to be necessary. Its inclusion in Nat is a matter for further debate and experimentation. Notice that with

permanent assignments, there is no possibility of mis-addressed packets, and so the two-sided technique has no effect there.

3.3 Running Routing Algorithms Across Nat

Of course, in order for Nat to be transparent to the border routers of the backbone and the stub, the border routers must believe that they are exchanging routing information with each other in the usual way. In other words, the stub border router must think that it is sending routing information about its internal (local) addresses to the backbone border router. Nat must intercept the routing information from the stub border router and replace the local address information with address information reflecting its global pool. However, global information that Nat receives from the stub border router must be passed through unchanged. Routing information from the backbone border router to Nat should always be passed through unchanged.

3.4 Multiple Nat Boxes and DNS Servers

All of the previous descriptions assume that each stub has just one Nat box and one DNS server. However, each stub/backbone entry/exit point needs to have a Nat box. In many, if not most cases, an IP packet can potentially travel through more than one border router. For this to work in the context of Nat, every Nat box that might potentially handle the packets of a given connection must know the assignment.

In addition, any given host may have multiple DNS servers (primary and backup, for instance). In what follows, we describe the mechanisms necessary to make multiple Nat boxes and DNS servers work.

Need for Nat Cliques

The Nat boxes of a stub are partitioned into one or more possibly overlapping groups, each with one or more Nat boxes. We call these groups Nat cliques. A Nat clique is a group of Nats such that a packet addressed with an assignment from one of the Nats in the clique can potentially be routed through any of the Nats in the clique. Therefore, the formation of a Nat clique depends on both intra-domain and inter-domain routing, primarily inter-domain.

There can be many reasons why such cliques form. For instance, assume that a stub has two attachments each to both NSFNET and MILNET. Assume also that the addresses assigned to the Nat boxes corresponding to each backbone are hierarchically formed to imply routing through that backbone. Therefore, packets assigned by the NSFNET-attached Nats will go only through NSFNET, and packets

assigned by the MILNET-attached Nats will go only through MILNET. In this case, there are two Nat cliques.

The Nats in a Nat clique are divided into two types—those that can assign and receive addresses for the clique, and those that only receive addresses for the clique. We call these clique-assigning Nats and clique-receiving Nats.⁹ Although it is not absolutely necessary, every Nat should be a clique-assigning Nat for at least some clique (for instance, for robustness in case all other Nats in a clique fail).

The reason for having two types of Nats is as follows. Consider the previous example, but modify it so that it is possible to alternate route packets with the MILNET-derived assignments through NSFNET, but it is not possible to alternate route packets with the NSFNET-derived assignments through MILNET. In this case, the NSFNET Nats must be aware of the MILNET assignments, in case such packets are alternate routed through them, but the MILNET Nats do not need to be aware of the NSFNET Nat assignments. There are therefore two cliques. The “NSFNET” clique has just the two NSFNET-attached Nats, and both are clique-assigning Nats. The “MILNET” clique has all four Nats, but the MILNET-attached Nats are clique-assigning Nats, and the NSFNET-attached Nats are clique-receiving Nats.

Operation of Nat Cliques and DNS Cliques

Nat cliques and DNS cliques require the following configuration information. The DNS server contains a list for each Nat clique (in the case where interception is not being used). Within each list is the IP address of the clique-assigning Nats in the clique. Each Nat contains a list for each Nat clique for which it is a clique-assigning member. This list contains the IP address of every Nat in the clique (both clique-assigning and clique-receiving). Clique-assigning Nats must also contain a list for each DNS clique (in the case where interception is not being used). A DNS clique is a group of DNS servers that can potentially answer a query for the same hosts. Although it is only necessary as a configuration-correctness mechanism, each Nat may contain a list for each Nat clique for which it is a clique-receiving member. This list contains the IP address for every clique-assigning member of the clique. All of the configuration information must be stored in non-volatile memory (or be learnable upon booting).

Every Nat box has one or more unique pools of global addresses from which it can make assignments. By not having Nat boxes share global addresses, we eliminate the need for coordination of address assignment

9. If Nat intercepts DNS packets, then there can be only one clique, and all members of the clique must be clique-assigning Nats. This is because a DNS packet can be routed through any Nat box, and so all Nat box must be prepared to make assignments.

where one Nat box has to check with others to make sure that a particular assignment can be made. This has the negative effect of requiring a larger pool in each Nat box than what otherwise would be necessary (to insure that any Nat box doesn't run out of addresses to assign). However, the increase is not that much (certainly not linear with the number of Nat boxes), since assignments can be spread over the Nats in a clique.

The reason that a Nat box may have multiple pools is because a backbone may assign multiple address prefixes to a single stub for the purposes of policy routing. For instance, assume that a regional backbone is attached to both MILNET and NSFNET. If the regional network maintains two address prefixes, and advertises one of them to MILNET and the other to NSFNET, then packets with the MILNET-advertised address will be routed through MILNET and vice versa. By receiving multiple addresses in a query, a host (or user) has the power to choose the backbone network [Ts1].

When a DNS server receives a DNS query, it sends a Nat-assignment query to one clique-assigning Nat in each Nat clique. It should round-robin the queries among the clique-assigning Nats in each clique to evenly spread the assignment load. The Nat receiving the query makes an assignment, returns the assignment to the DNS server, sends an assignment notification message to all DNS servers that are members of the DNS cliques that the requesting DNS server belongs to, and sends an assignment notification message to all of the Nats in its cliques indicating the new assignment. On rare occasions, the Nat receiving the query may have no addresses available for assignment. In this case, Nat returns a NULL assignment, and the DNS may either query another Nat box in the clique, or return a failure in its DNS response.

When the assignments are returned to the DNS servers, they have expiration times associated with them. (If interception is used, then Nat sets the TTL in the DNS response itself.) The assigning Nat boxes must not reassign the address until the expiration time has elapsed. The DNS servers must not set the TTL (Time-to-Live) field in the DNS response to longer than the shortest expiration time. If the DNS servers choose to cache the assignment, they must remove the cache entry by the shortest expiration time. Note that it is not necessary for the DNS servers to cache the entry, because if another query for the same host comes, the DNS server will query Nat boxes and receive the same assignments.

When Nat boxes receive assignment notifications, they keep the assignments until notified otherwise. This will occur when the assigning Nats reassigns the addresses.

Nat assignment notifications must be reliable, because there is no refreshing (or timing out) of assignments by receiving Nats. Therefore, assignment notification messages must be acknowledged, and resent if no

acknowledgment is received. Of course, if a receiving Nat has crashed, then no acknowledgment can be sent. Therefore, Nat boxes must be able to mark other Nat boxes as down after a number of attempted assignment notifications. Also, when Nat boots (comes up after crashing), it must contact all assigning Nats in its cliques and receive all current assignments. This must also happen if Nats in a clique have been partitioned from each other, and the partition heals. Note that each Nat must have enough memory to hold all of the assignments of all of the Nats in all of their cliques.

Private Networks that Span Backbones

In many cases, a private network (such as a corporate network) will be spread over different locations and will use a public backbone for communications between those locations. In this case, it is not desirable to do address translation, both because large numbers of hosts may want to communicate across the backbone, thus requiring large global address pools, and because there will be more applications that depend on configured addresses, as opposed to going to a name server. We call such a private network a *backbone-partitioned stub*.

Backbone-partitioned stubs should behave as though they were a non-partitioned stub. That is, the routers in all partitions should maintain routes to the local address spaces of all partitions. Of course, the (public) backbones do not maintain routes to any local addresses. Therefore, the border routers must tunnel through the backbones using encapsulation. To do this, each Nat box will set aside one global address from the pool for tunneling. When a Nat box *x* in stub partition *X* wishes to deliver a packet to stub partition *Y*, it will encapsulate the packet in an IP header with a destination address from the pool of Nat box *y* that has been reserved for encapsulation. Then Nat box *y* receives a packet with that destination address, in decapsulates the IP header and routes the packet internally.

3.5 Header Manipulations

In addition to modifying the IP address, Nat must modify the IP checksum, the TCP checksum, places in ICMP and FTP where the IP address appears, and perhaps other places where the IP address appears¹⁰.

The checksum modifications to IP and TCP are simple and efficient. Since both use a one's complement sum, it is sufficient to calculate the arithmetic difference between the before-translation and after-translation addresses and add this to the checksum. The only tricky part is determining whether the

10. The author knows of no other such places off hand, but there are undoubtedly some. Hopefully, most such applications will be discovered during experimentation with Nat.

addition resulted in a wrap-around (in either the positive or negative direction) of the checksum. If so, 1 must be added or subtracted to satisfy the one's complement arithmetic. Sample code (in C) for this is as follows:

```
int16 nat_newChk (checksum, oldaddr, newaddr)
    int16 checksum;
    int32 oldaddr, newaddr;
{
    int16 newCheck, oldCheck;
    int32 chk32, diff32;
    int16 crossing;
    int32 carry1, carry2;

    /* diff32 could be pre-calculated when assignment is made */
    diff32 = (((newaddr >> 16) & 0x0000ffff) + (newaddr & 0x0000ffff)) -
        (((oldaddr >> 16) & 0x0000ffff) + (oldaddr & 0x0000ffff));

    checksum = ~checksum;
    chk32 = (0x0000ffff & checksum);
    chk32 += 0x00020000;
    chk32 += diff32;
    crossing = (chk32 >> 16) - 2;

    chk32 += crossing;
    checksum = 0xffff & (int16)chk32;

    return (~checksum);
}
```

The arguments to the File Transfer Protocol (FTP) PORT command include an IP address (in ASCII!). If the IP address in the PORT command is the same as that of the host sending the PORT command, then Nat must substitute the (local) IP address in the FTP PORT command with the (global) assigned IP address. Because the address is encoded in ASCII, this may result in a change in the size of the packet (for instance, 10.18.177.42 is 12 ASCII characters, while 193.46.228.137 is 14 ASCII characters). If the packet size is changed in transit, then the subsequent TCP sequence numbers (which are in units of bytes) will be wrong, and TCP will fail.

In some cases, it may be possible for Nat to choose a global IP address that has the same number of ASCII characters as the local IP address. It is possible, however, for the character size of the local IP address to be smaller (or larger) than the smallest (or largest) possible IP address from the pool. In this case, FTP will fail. In general, in order to run FTP outside of a stub, it will be necessary to either limit outside FTP to a few internally widely available hosts, or set up an FTP application gateway.

If the IP address in the PORT command is different from that of the host sending the PORT command, but the IP address is local to the stub domain, then Nat can create an assignment for the IP address and substitute that. Since the address is encoded in ASCII, the TCP checksum cannot as easily be incrementally recalculated, and should therefore be recalculated from scratch. If the IP address in the PORT command is not from the local stub, then it should not be modified. Of course, if the FTP session is encrypted, the PORT command will fail.

If an ICMP message is passed through Nat, it may require two address modifications and three checksum modifications. This is because most ICMP messages contain part of the original IP packet in the body. Therefore, for Nat to be completely transparent to the host, the IP address of the IP header embedded in the data part of the ICMP packet must be modified, the checksum field of the same IP header must correspondingly be modified, and the ICMP header checksum must be modified to reflect the changes to the IP header and checksum in the ICMP body. Of course, the normal IP header must also be modified as already described.

It is not entirely clear that the IP header information in the ICMP part of the body really need be modified. This depends on whether or not there is really any host code that looks at this IP header information¹¹. It may in fact be useful to not translate, so as to provide the exact header seen by the router or host that issued the ICMP message, which may aid in debugging. In any event, no modifications are needed for the Echo and Timestamp messages, and Nat should never need to handle a Redirect message.

3.6 Other Aspects of Nat

Global Routing and Addressing Issues

Over the short term, Nat provides scaling benefits by allowing for subnetting of stubs by backbone networks. In doing this, we essentially add a level of hierarchy to IP routing. We also introduce the coupling of route to address that the OSI community is now having to face. Namely, if an IP address is handed out by a backbone, and that backbone advertises that address as reachable through it, then routes will naturally go through that backbone. If an alternate route through another backbone is desired (for instance, because the primary route failed), that route may not be available.

11. In the theoretical worst case, an ICMP message could be sent concerning an FTP packet that contained a PORT command. In this case, modifications would be required to the PORT command and the TCP checksum, in addition to the fields already mentioned. In practice, this seems unnecessary.

Viewed another way, this coupling of route to address can actually be a feature rather than a bug. If a host or user wishes to route through one backbone vs. another, it can manipulate the choice by choosing the appropriate address. This would work as follows. When the DNS server queries the Nat boxes for assignments, it may get back multiple answers, one from each Nat clique, and possibly multiple assignments from a single Nat clique. These multiple answers essentially reflect reachability of the stub through multiple backbones. When the DNS server then returns the queries, the source host can choose the appropriate one.

The use of multiple addresses as a means of policy routing and scaling are discussed extensively in [Ts1]. The main point here is that the extra layer(s) of IP address hierarchy resulting from Nat make it possible to take advantage of multiple addresses.

Dynamic Allocation of Nat Pool

The size of the pool of addresses needed by a Nat box varies over time. At certain times more addresses are needed than at others. If the Nat pools can be dynamically assigned to Nat boxes from a larger pool, then the benefits of statistical sharing can be realized. Each Nat box could keep a pool large enough to handle most of its needs, but the Nat box could dynamically request more addresses from its backbone when necessary.

This could be done using the two-level kampai algorithm described in [Ts2]. That algorithm is for the purpose of assigning subnet numbers. Its main advantage is that it allows subnet numbers to be assigned efficiently without requiring advance knowledge of the size (in terms of number of hosts) and number of subnets. It does this by removing a bit from the mask when more space is needed in a subnet. This doubles the subnet's space. Since this algorithm can be automated, it may be possible for Nat boxes to request and return address space in increments of powers of two.

Applications with IP-address Content

Any application that carries (and uses) the IP address inside the application will not work through Nat unless Nat knows of such instances and does the appropriate translation. It is not possible or even necessarily desirable for Nat to know of all such applications. And, if encryption is used then it is impossible for Nat to make the translation.

It may be possible for such systems to avoid using Nat, if the hosts in which they run are assigned global addresses. Whether or not this can work depends on the capability of the intra-domain routing algorithm and the internal topology. This is because the global address must be advertised in the intra-domain routing algorithm. With a low-feature routing algorithm like RIP, which does not use masks, the host requires its own class C address space. This address must be advertised externally as well as internally (thus hurting global scaling). With a high-feature routing algorithm like OSPF, which does use masks, the host address can be passed around individually, and can come from the Nat pool.

Privacy, Security, and Debugging Considerations

Unfortunately, Nat reduces the number of options for providing security. With Nat, nothing that carries an IP address or information derived from an IP address (such as the TCP-header checksum) can be encrypted. While most application-level encryption should be ok, this prevents encryption of the TCP header.

On the other hand, Nat itself can be seen as providing a kind of privacy mechanism. This comes from the fact that machines on the backbone cannot monitor which hosts are sending and receiving traffic (assuming of course that the application data is encrypted).

The same characteristic that enhances privacy potentially makes debugging harder (including tracking security violations) more difficult. If a host is abusing the Internet in some way (such as trying to attack another machine or sending large amounts of junk mail or something) it is more difficult to pinpoint the source of the trouble because the IP address of the host is hidden.

4.0 Conclusions and Status

Nat may be a good short term solution to the address depletion and scaling problems. This is because it requires no changes to existing hosts and routers, or only changes to DNS, and can be installed incrementally. Nat has several negative characteristics that make it inappropriate as a long term solution, and may make it inappropriate even as a short term solution. Only implementation and experimentation will determine its appropriateness.

The negative characteristics are:

1. Not all applications can run over Nat. The most serious case is with FTP, where the ASCII size of the substituted IP address sometimes cannot be made the same as the original IP address (thus breaking TCP sequence numbering). While it is possible to run FTP without the PORT command, this requires changes in current implementations. In addition, it won't work with encryption of the TCP header, or encryption of any upper layer protocols that encode IP addresses for the sake of determining the return IP address.
2. It requires a sparse end-to-end traffic matrix. Otherwise, the Nat pools will be large, thus using up many addresses. While the expectation is that end-to-end traffic matrices are indeed sparse, experience with Nat will determine whether or not they are. In any event, future applications may require a rich traffic matrix (for instance, distributed resource discovery), thus making long-term use of Nat unattractive.
3. It may significantly increase the load on DNS. This is because DNS server will not be able to cache responses for as long. Currently, DNS represents a significant proportion of Internet traffic. Since current caching efficiency for DNS is not currently known, the DNS traffic increase may be light, or it may be heavy.
4. It increases the probability of mis-addressing.
5. It won't work with hosts that don't query DNS (except for permanent assignments).
6. It hides the identity of hosts. While this has the benefit of privacy, it is generally a negative effect.

4.1 Current Implementation

An experimental prototype of Nat is being implemented on public domain KA9Q TCP/IP software [Ka]. The prototype currently does permanent assignments (administered from the system monitor), but does not do dynamic assignment in concert with DNS queries. The prototype has demonstrated that IP addresses can be translated transparently to hosts within the limitations described in this paper.

But we have not gathered any data having to do with the dynamics of address assignment—assignment strategies, appropriate pool sizes, and mis-addressed packets. Since this is the more interesting aspect of Nat experimentation, much work remains to be done.

Acknowledgments

The authors would like to acknowledge Van Jacobsen, who initially expressed the concept of address reuse in the context of IP, and Dave Clark for his review of this work.

REFERENCES

- [Ch] Chiappa, N., "The IP Addressing Issue", IETF Internet Draft, draft-chiappa-ipaddressing-00.txt, anonymous FTP from nnsc.nsf.net, March, 1991.
- [Ka] Karn, P., "KA9Q", anonymous FTP from ucsd.edu (hamradio/packet/ka9q/docs).
- [Lo] Lottor, M., "Domain Administrators Operations Guide", RFC-1033, USC/Information Sciences Institute, November 1987.
- [Mo] Mockapetris, P.V., "Domain names - implementation and specification", RFC-1035, USC/Information Sciences Institute, November 1987.
- [Pa] Paxton, V., "Measurements and Models of Wide Area TCP Conversations", LBL-30840, Lawrence Berkeley Laboratory, Berkeley, California, May, 1991.
- [Ts1] Tsuchiya, P.F., "Robust and Efficient Policy Routing using Multiple Hierarchical Addresses", ACM SIGCOMM '91, Zurich, Switzerland, September 1991.
- [Ts2] Tsuchiya, P.F., "On the Assignment of Subnet Numbers", RFC 1219, USC/Information Sciences Institute, April 1991.